

## ABSTRACT

A method and apparatus for transforming a web page that contains main content and auxiliary data. The web page is converted into a string containing multiple first values and multiple second values. The first values correspond to formatting code segments within the web page and the second values correspond to text segments within the web page. Further, a low-pass filter is applied to the string containing multiple first values and multiple second values, and the output of the low-pass filter is used to determine the location of the main content within the web page.

Downloaded from ascelibrary.org by University of California, San Diego on 06/01/15. Copyright ASCE, For All Rights Reserved, No part of this document may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or by any information storage or retrieval system, without permission in writing from ASCE.